

Enabling Big Geoscience Data Analytics with a Cloud-based, MapReduce-enabled and Service-oriented Workflow Framework

- An Example with Climate Model Sensitivity Analysis

Zhenlong Li¹, Chaowei Yang¹, Baoxuan Jin^{1,2}, Manzhu Yu¹, Kai Liu¹, Min Sun¹,
Matthew Zhan^{1,3}

1. NSF Spatiotemporal Innovation Center, George Mason University, Fairfax, VA, USA

2. Yunnan Provincial Geomatics Center, Yunnan Bureau of Surveying, Mapping, and
GeoInformation, Kunming, Yunnan, China

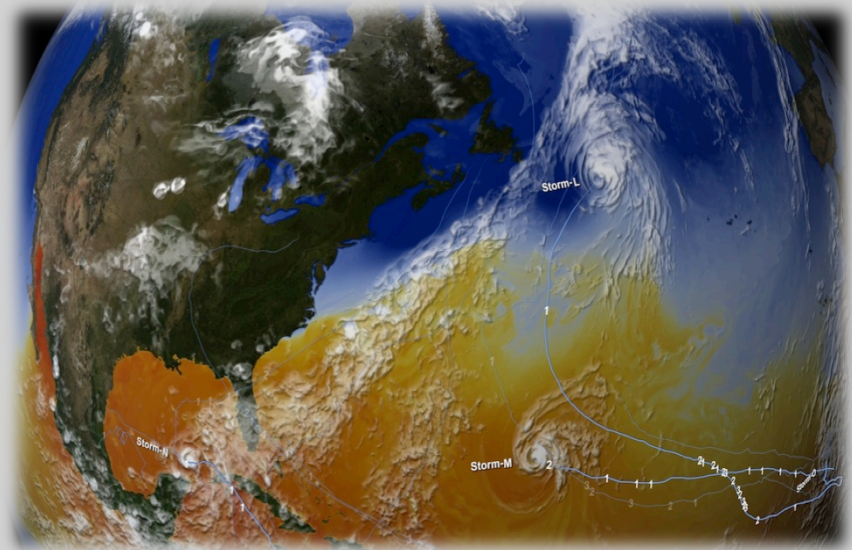
3. Department of Computer Science, University of Texas – Austin, Austin, Texas, USA

Agenda

1. Climate Model Sensitivity Analysis
2. Workflow Framework
3. Result and Demonstration
4. Conclusion

Climate Model Sensitivity Analysis

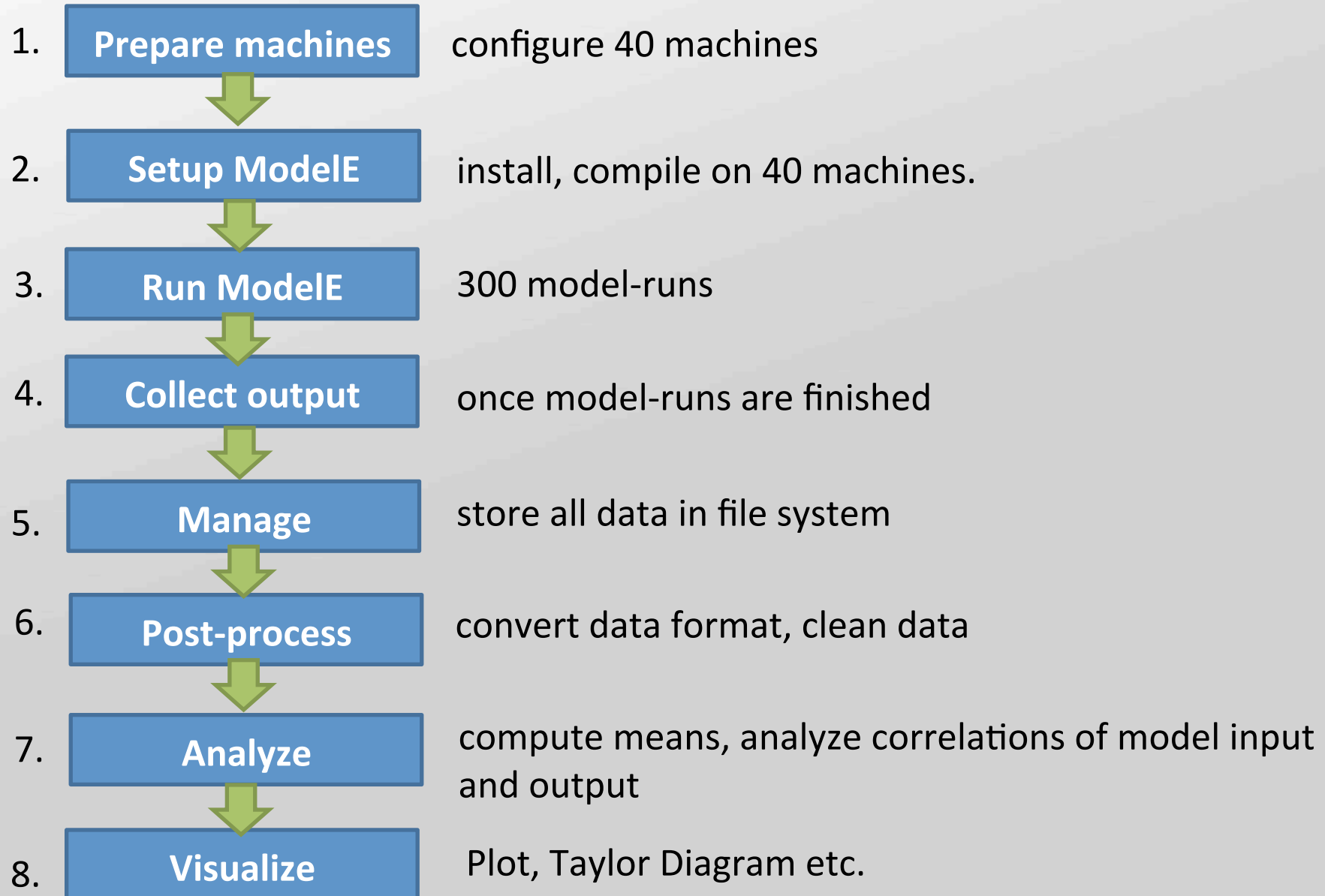
- An essential procedure in climate modeling
- Understand relationships between model input and output
- Run a model many times by sweeping the model input space



An Example

- ModelE: a global climate model developed by NASA Goddard Institute of Space Studies (GISS)
- Test ModelE sensitivity with perturbed physics ensemble(PPE) experiment
- Run ModelE 300 times with different parameter values(PPE-300)
- Conduct PPE-300 many times to investigate how different parameter combinations impact simulated climate

Workflow for the PPE-300



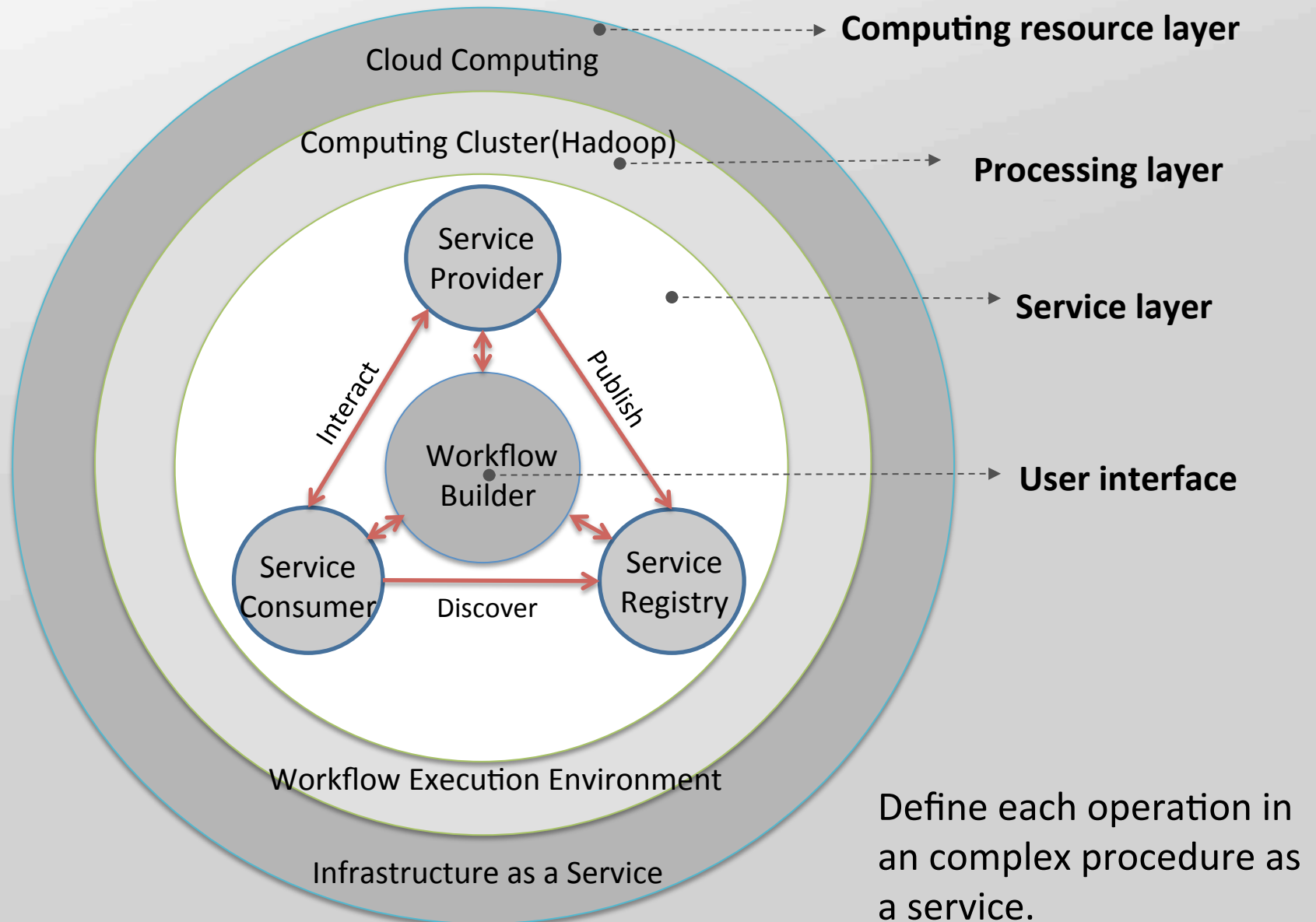
Took two months to finish PPE-300!

A Ph.D. student from GMU collaborated with a NASA scientist

Challenges

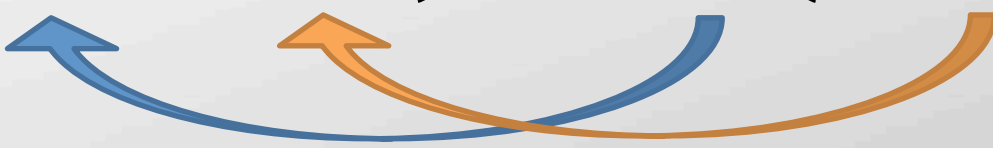
- Computing Intensity
 - 300 model-runs → 1500 days (4 years)
- Data Intensity
 - One PPE-300 experiment → 2.5 terabytes climate data
- Procedure Complexity
 - Multiple steps, different tools, libraries, external services

Cloud-based, MapReduce-enabled and Service-oriented Workflow Framework



Unified Service Model

Service(*DataIn*, *ParaIn*) \rightarrow *execute*⁺ (*DataOut*, *ParaOut*)



Four service types

Data Service

- Fetch data from external sources(FTP, OpenDAP)
- Publish data

Process/Analytic Service

- Parallel processing(MapReduce)
- R script, Python program ...

Model Service

- Enabled by MaaS(Model as a Service)

Infrastructure Service

- Start a virtual machine
- Start a computing cluster

Standard Service Metadata

The image shows a text editor window displaying XML service metadata. The XML is for a service named 'ProvisionVM' and is annotated with three red boxes and labels:

- Service description:** This box highlights the `<description>`, `<category>`, `<provider name>`, and `<contact>` elements.


```

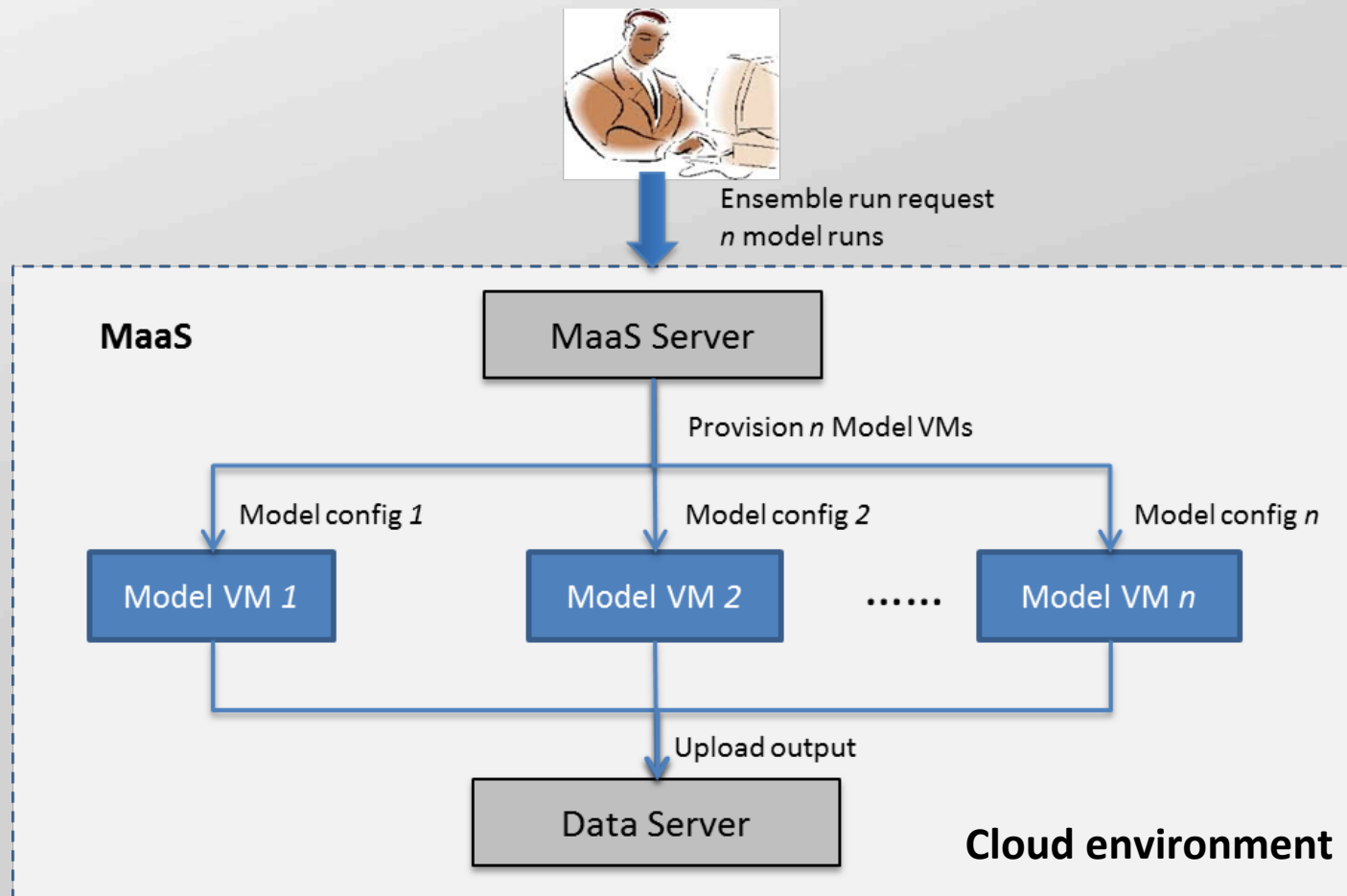
      <?xml version="1.0"?>
      <!--A standard xml for describing the service metadata for the workflow platform -->
      - <service name="ProvisionVM" xmlns="uri:cisc:service:0.1">
        <description>Provision a virtual machine based on user's configuration </description>
        <category>Cloud</category>
        - <provider name="CISC">
          <site>http://cisc.gmu.edu</site>
          - <contact>
            <individual-name>Zhenlong Li</individual-name>
            <email>zhenlong.gis@gmail.com</email>
            <phone>703-909-3676</phone>
          </contact>
        </provider>
      
```
- Service entry point:** This box highlights the `<entry-point type>` element.


```

      - <entry-point type="JAVA">
        <main-class>gmucisc.service.ProvisionVM</main-class>
        <jar-location>runnable jar file in a web accessible location</jar-location>
      </entry-point>
      <!--the order for the inputs matters, for the outputs doesn't matter-->
      - <interface>
        - <input name="VM type">
          <description>Virtual machine types, following the standard of Amazon EC2, value could be m1.small, c1.medium, c1.xlarge</description>
          <default-value>m1.small</default-value>
        </input>
        - <output name="VM instance id">
          <description>instance id for the virtual machine, can be used to terminate the virtual machine</description>
          <variable-id>VM_ID</variable-id>
        </output>
        - <output name="VM public ip">
          <description>public IP address for this vm</description>
          <variable-id>VM_IP</variable-id>
        </output>
      </interface>
    </service>
    
```
- Service interface: input/output:** This box highlights the `<input>` and `<output>` elements within the `<interface>` block.

Model Service(MaaS)

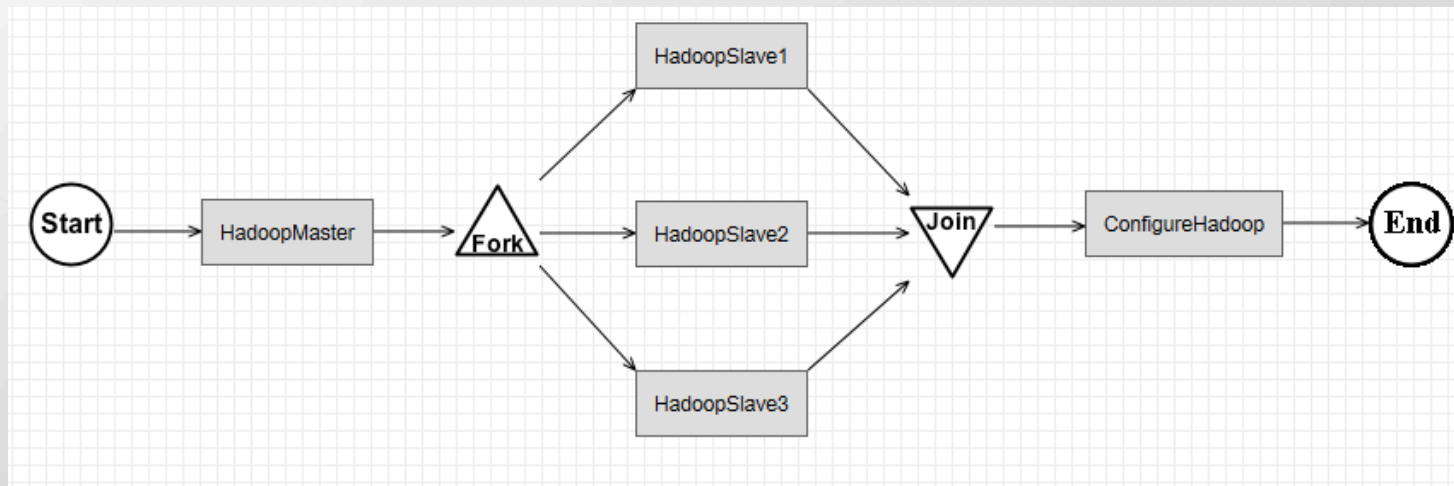
Hundreds of virtual machines with model environment will be started in minutes to run models in parallel



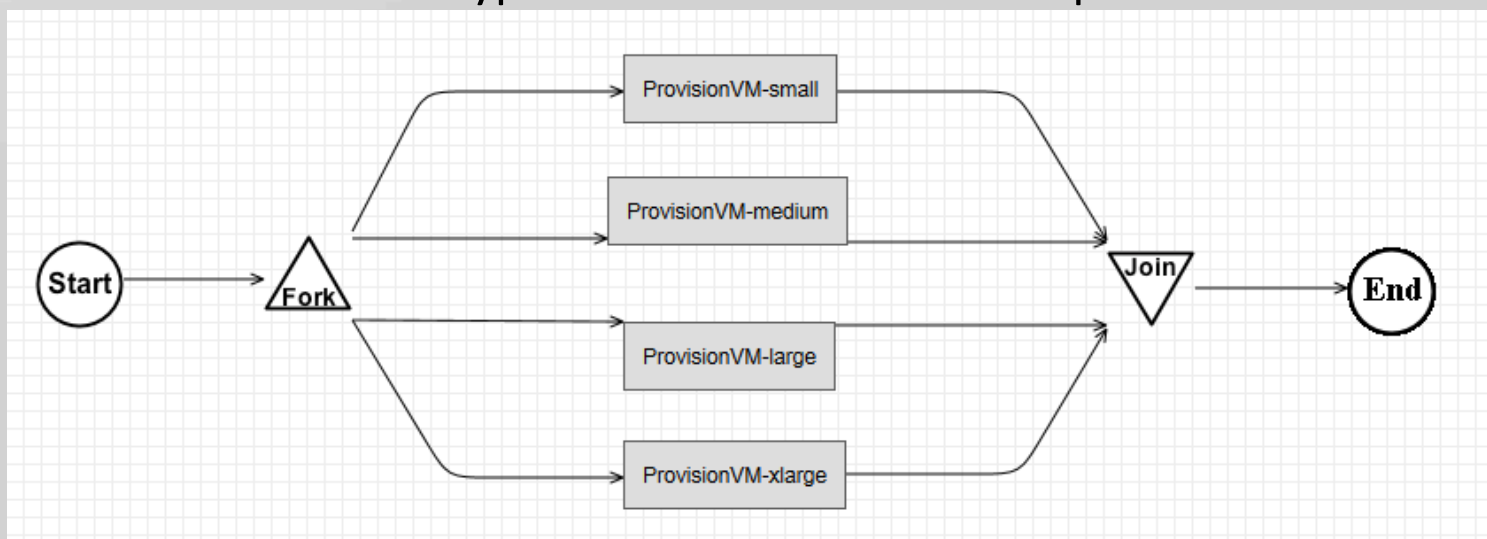
Li, Z., Yang, C., Huang, Q., Liu, K., Sun, M., & Xia, J. (2014). Building Model as a Service to support geosciences. *Computers, Environment and Urban Systems*.

Infrastructure Service

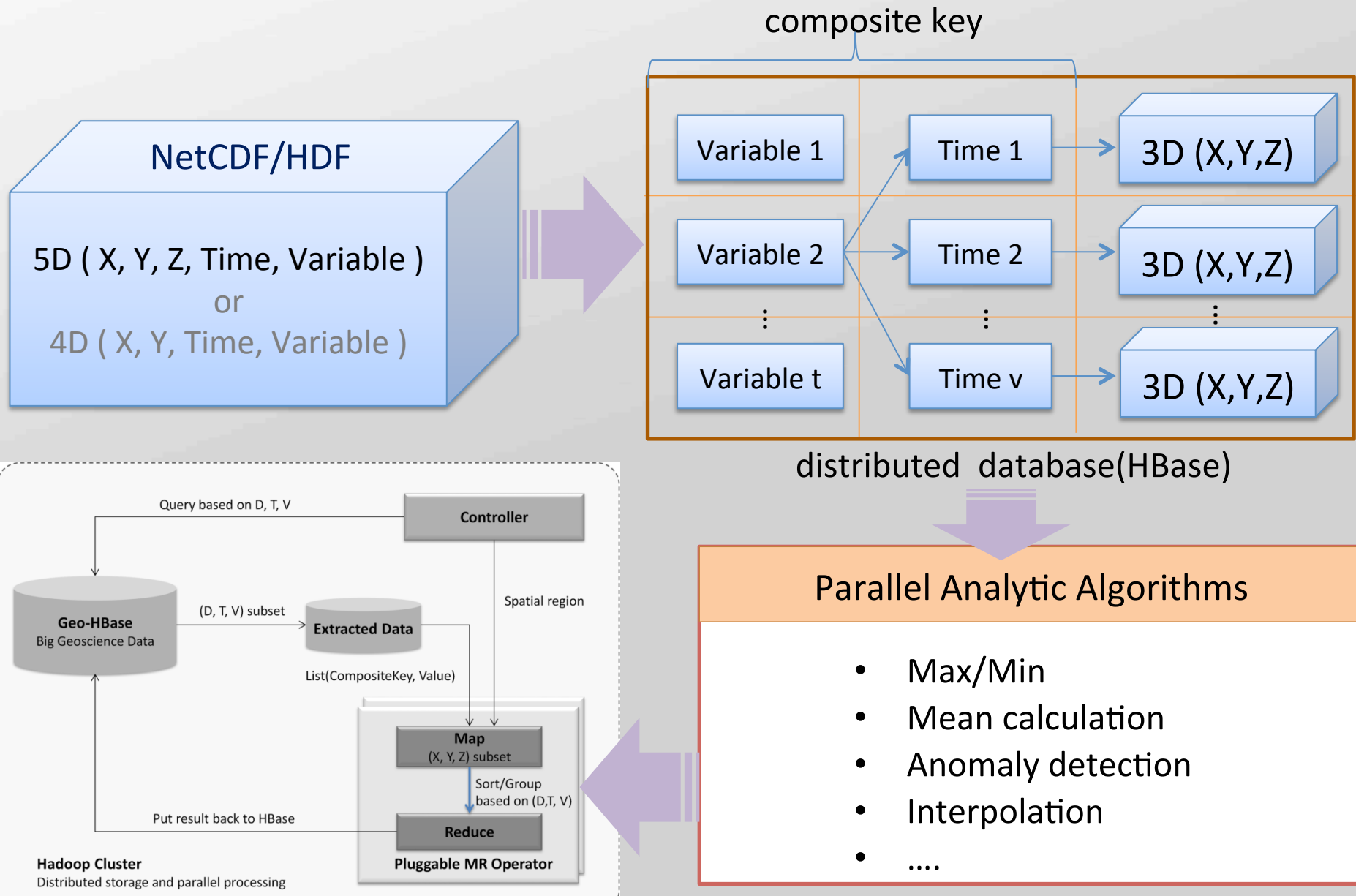
Start a 3-node computing cluster (Hadoop)



Start four different types of virtual machines in parallel



Big Climate Data Processing with MapReduce

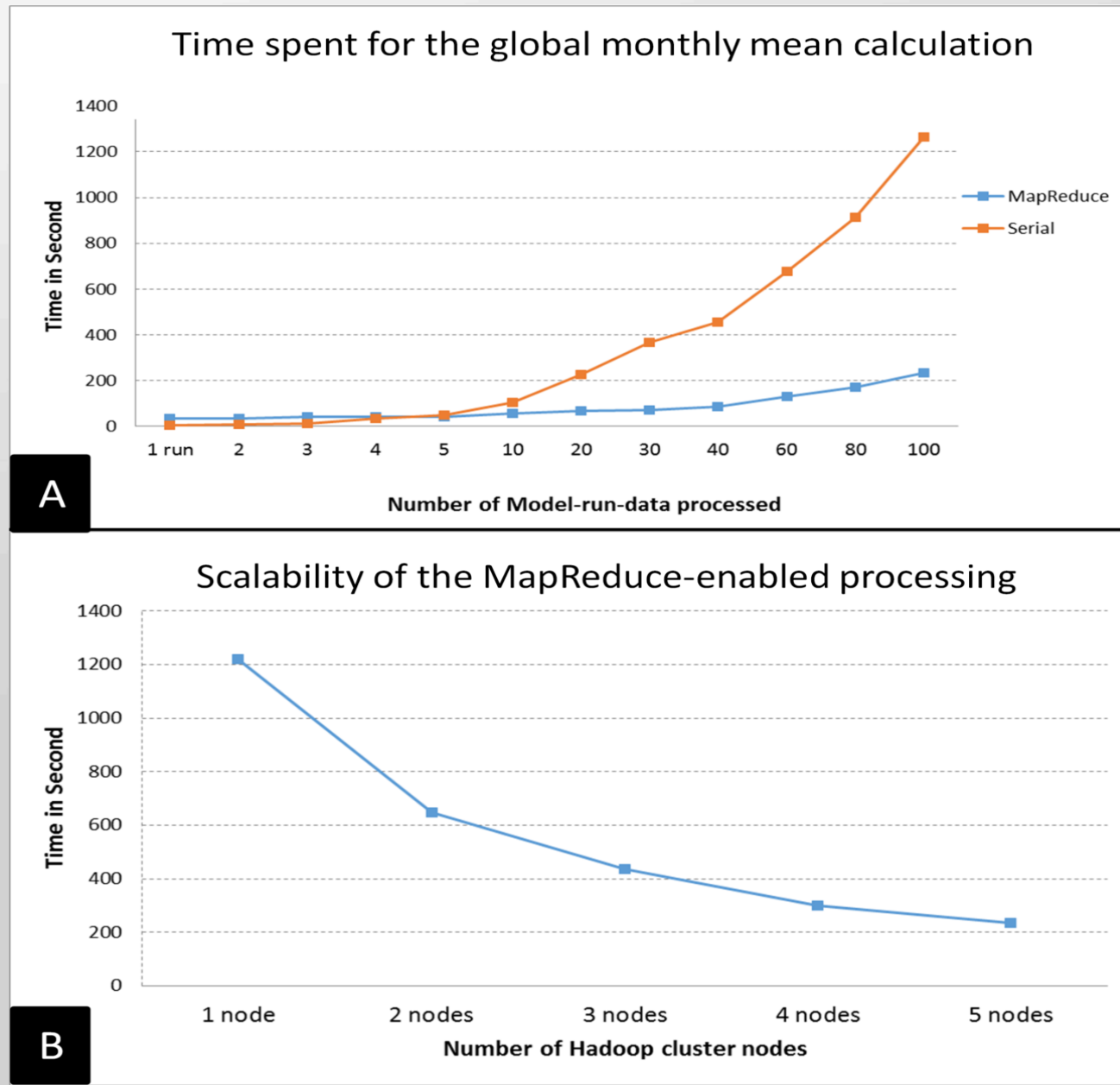


Result: ensemble model run in the cloud

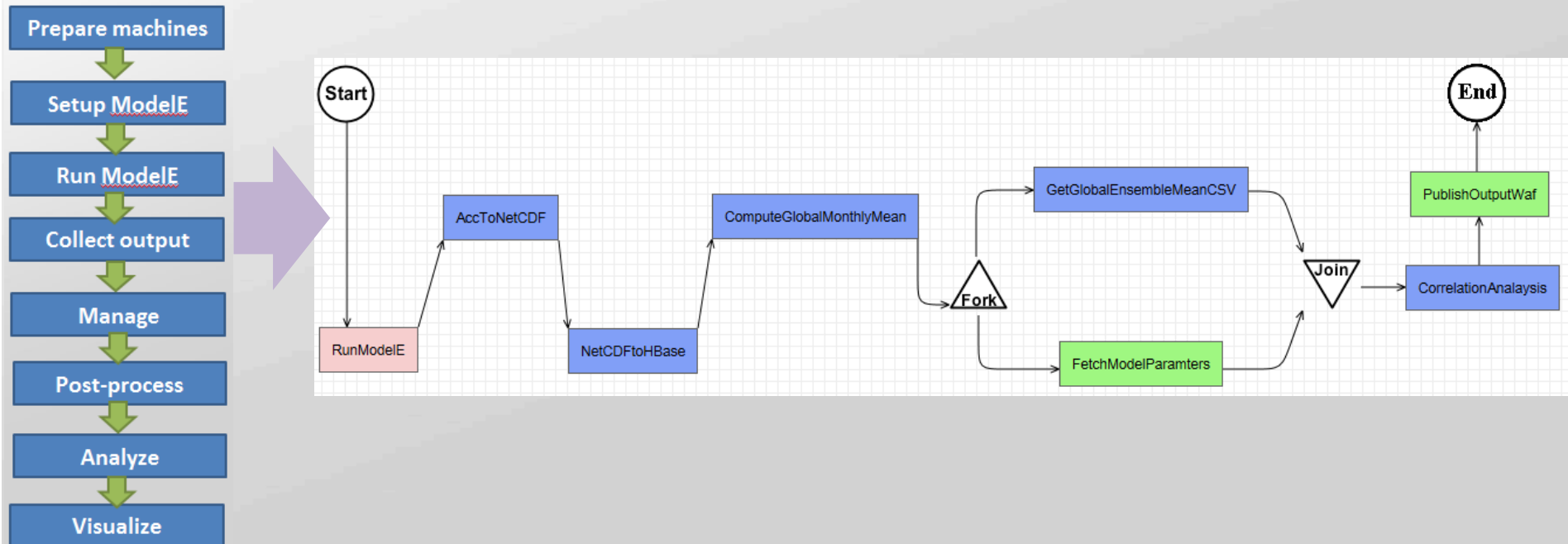
300 model-runs with/without MaaS using the same 40 machines

Comparing Item	Without MaaS	With MaaS
Model setup time	~9 days	~10 minutes
300 model-run time	~38 days	~5 days
Resource(CPU) utilization rate	12%	93%

Result: mapreduce-enabled processing



Result: runnable workflow for PPE-300



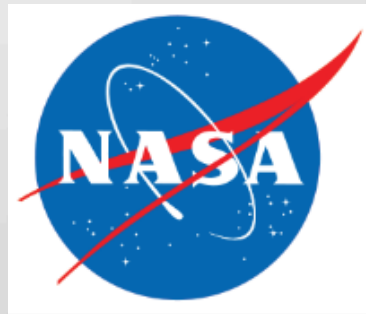
Workflow Prototype Demo

- A video showing the PPE-300 experiment.
- A live demo showing the MERRA data analysis(provision cluster, download data, analyze and visualize).

Conclusion

- **Cloud-based for computing intensity**
 - Dynamically provide computing resources during workflow execution
- **MapReduce-enabled for data intensity**
 - Analyze big climate data in parallel with a highly scalable framework
- **Service-oriented for procedure complexity**
 - Unified Service Model transforms complex experiment into executable workflows

Acknowledgements



Q&A

Thank you!